# Problems and Paradoxes in High-dimensional Data

Michael Komodromos

February 7, 2022

# Outline

# Sources of high-dimensional data

# Sources of high-dimensional data

- **Biotech**: DNA microarrays, proteomics, transcriptomics etc.

# Sources of high-dimensional data

- **Biotech**: DNA microarrays, proteomics, transcriptomics etc.

- **Images / video**: medical, astrophysics, surveillance

# Sources of high-dimensional data

- **Biotech**: DNA microarrays, proteomics, transcriptomics etc.

- **Images / video**: medical, astrophysics, surveillance

- **Consumer preferences**: books, movie, music recommendation

# Sources of high-dimensional data

- **Biotech**: DNA microarrays, proteomics, transcriptomics etc.

- **Images / video**: medical, astrophysics, surveillance

- **Consumer preferences**: books, movie, music recommendation

- etc.

# Examples

240p image $\sim 100,000$ pixels, 3 color bandwidth $p \approx 300,000$

# Paradoxes

# Volume of an n-D ball

- What's the volume of a circle with radius $r$?

# Volume of an n-D ball

- What's the volume of a circle with radius $r$?

- $\pi r^2$

# Volume of an n-D ball

- What's the volume of a circle with radius $r$?
- $\pi r^2$

- What's the volume of a ball with radius $r$?

# Volume of an n-D ball

- What's the volume of a circle with radius $r$?

- $\pi r^2$

- What's the volume of a ball with radius $r$?

- $\frac{4}{3}\pi r^3$

# Volume of an n-D ball

- What's the volume of a circle with radius $r$?

- $\pi r^2$

- What's the volume of a ball with radius $r$?

- $\frac{4}{3}\pi r^3$

Volume of $p$-dimensional ball with radius $r > 0$

$$V(p; r) = \frac{1}{\Gamma(p/2 + 1)}\pi^{p/2} r^p \tag{1}$$

Does the volume get bigger or smaller as we increase $p$?

Or more generally, how does the volume behave?

# The volume

# The volume

# Where is the mass concentrated?

We've noticed that the volume of a *p*-dimensional ball tends to 0 as *p* increases, then naturally we might ask:

Where is the mass concentrated?

# Where is the mass concentrated?

# Where is the mass concentrated?

Consider the volume in the crust, i.e. the volume

$$C(p; r) = V(p; r) - V(p; 0.99r) \qquad (2)$$

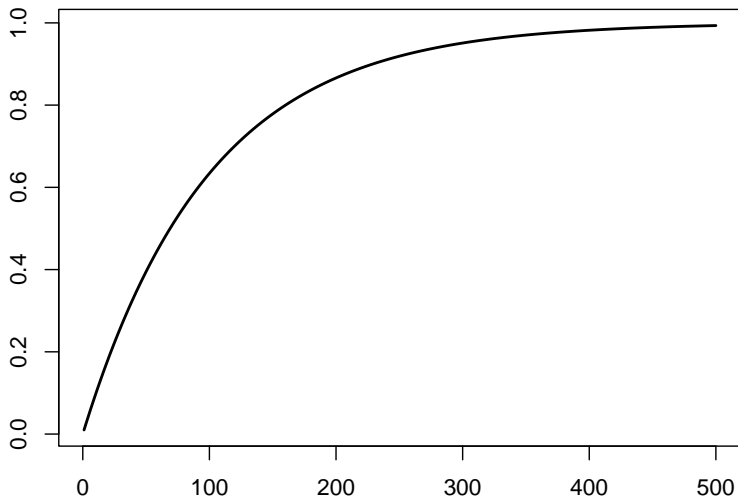# Where is the mass concentrated?

Consider the volume in the crust, i.e. the volume

$$C(p; r) = V(p; r) - V(p; 0.99r) \qquad (2)$$

## Where is the mass concentrated?

Consider the volume in the crust, i.e. the volume

$$C(p; r) = V(p; r) - V(p; 0.99r) \qquad (2)$$



Then the volume in the crust as a fraction of the total volume, i.e.

$$\frac{C(p; r)}{V(p; r)} = 1 - 0.99^p \qquad (3)$$

# Where is the mass concentrated?

# Where is the mass concentrated?

# So far...

# So far...

## Recap

- The volume of $p$-dimensional balls tends to 0

- The crust contains nearly all the mass

# So far...

## Recap

- The volume of $p$-dimensional balls tends to 0

- The crust contains nearly all the mass

The lesson in all this is, we must be careful with our **geometric intuition** of high-dimensional spaces!

# Paradoxes

# Paradoxes

There are many more counter-intuitive examples in high-dimensional spaces, to list a few

# Paradoxes

There are many more counter-intuitive examples in high-dimensional spaces, to list a few

- Most of the mass of a standard Gaussian is in the tail!

# Paradoxes

There are many more counter-intuitive examples in high-dimensional spaces, to list a few

- Most of the mass of a standard Gaussian is in the tail!
- Rare events may not actually be that rare

Problems

# Data get far apart fast

# Data get far apart fast

Suppose we have data uniformly distributed on a 2-D grid, 10-D grid, 100-D grid or 1000-D grid.

# Data get far apart fast

Suppose we have data uniformly distributed on a 2-D grid, 10-D grid, 100-D grid or 1000-D grid.

What's the average (euclidean) distance between points?
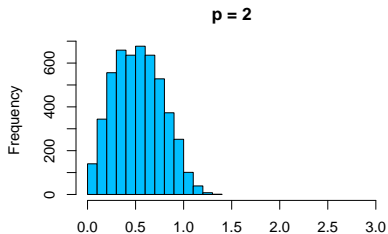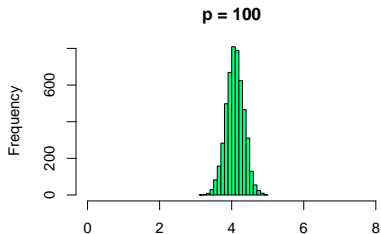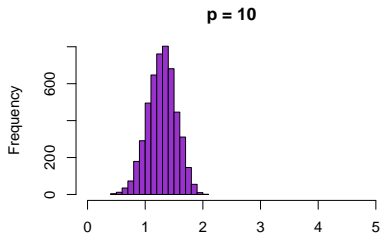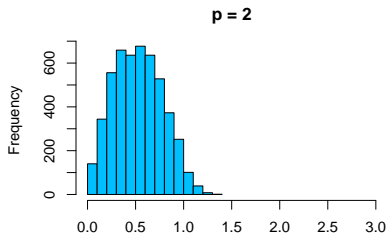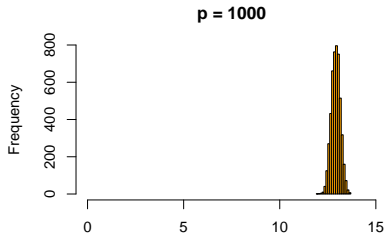
# 2D Case
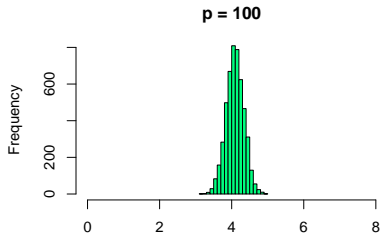
# 2D Case

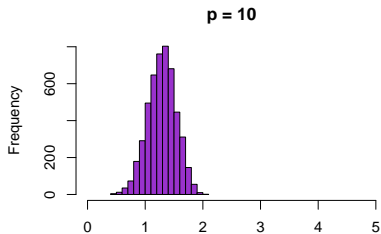# 2D Case
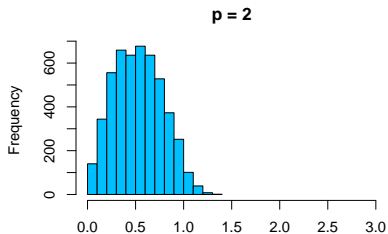
# 2D Case

# 2D Case

# What's the average distance?

# What's the average distance?

# What's the average distance?
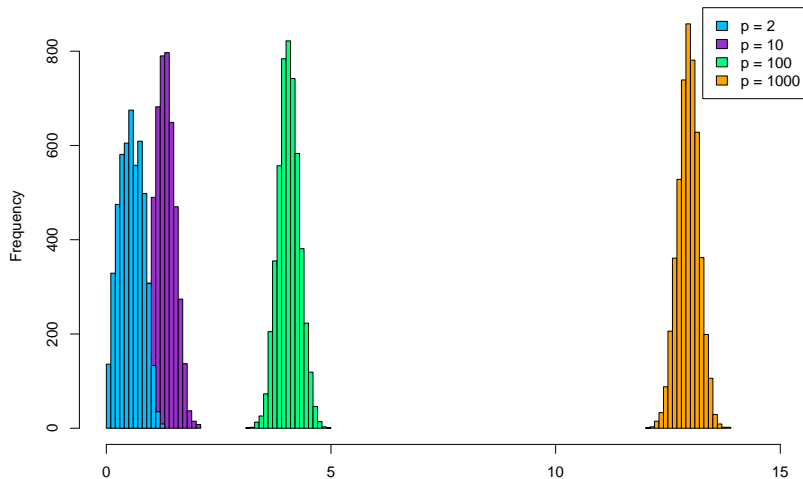
# What's the average distance?

# What's the average distance?

# What's the average distance?

# What's the average distance?

# What does this mean?

# What does this mean?

## Distances between points

- The minimal distance between two points increases

# What does this mean?

## Distances between points

- The minimal distance between two points increases

- All points are a similar distance from the others

# What does this mean?

## Distances between points

- The minimal distance between two points increases

- All points are a similar distance from the others

- The notion of nearest point vanishes

# Some more problems

# Some more problems

High-dimensional spaces are immense:

- Vast, data points can be isolated in the immensity

- Small changes add up fast. Many small fluctuations in different directions can produce a large global fluctuation

- Computation suffers

A practical example

# PCA

# PCA

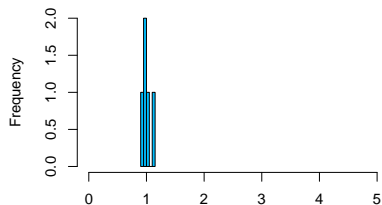PCA relies on the estimation of the covariance matrix $\Sigma$.

# PCA

PCA relies on the estimation of the covariance matrix $\Sigma$.

Let's suppose we have some some data generated iid from Normal$(0, I_p)$, where $I_p$ is the $p \times p$ identity.

# PCA
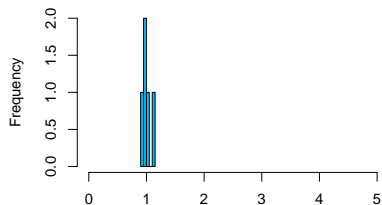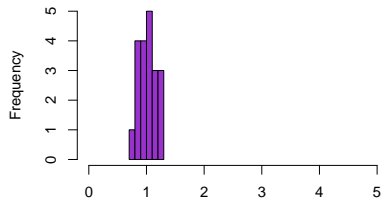
PCA relies on the estimation of the covariance matrix $\Sigma$.

Let's suppose we have some some data generated iid from Normal$(0, I_p)$, where $I_p$ is the $p \times p$ identity.

<span style="color:blue">How well is $\Sigma$ estimated as p grows?</span>
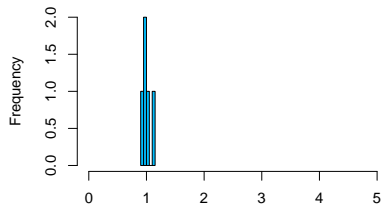
# PCA

# PCA

To see how well we're doing we're going to compare the eigenvalues of the estimated covariance matrices to that of the true covariance matrix $(I_p)$, we let $n = 1000$ and let $p = 5, 20, 250, 500$.
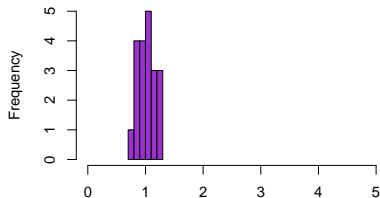
# PCA

# PCA

# PCA

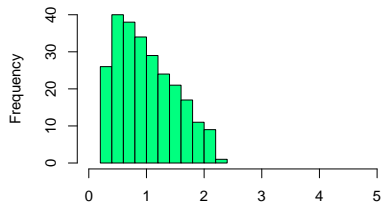# PCA
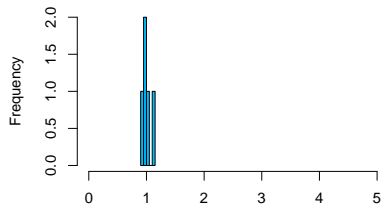
# PCA

# PCA

PCA, main message

# PCA

## PCA, main message

The empirical covariance is a very poor approximation of the covariance $I_p$ in this setting.

# PCA

### PCA, main message

The empirical covariance is a very poor approximation of the covariance $I_p$ in this setting.

In turn, downstream use of the covariance matrix (e.g. in PCA) can lead to spurious results.

# Silver lining

# Silver lining

In many settings data is often much more low-dimensional and not uniformly spread!

# Silver lining

In many settings data is often much more low-dimensional and not uniformly spread!

- Images have structures

# Silver lining

In many settings data is often much more low-dimensional and not uniformly spread!

- Images have structures

- Biological systems are strongly regulated

# Silver lining

In many settings data is often much more low-dimensional and not uniformly spread!

- Images have structures

- Biological systems are strongly regulated

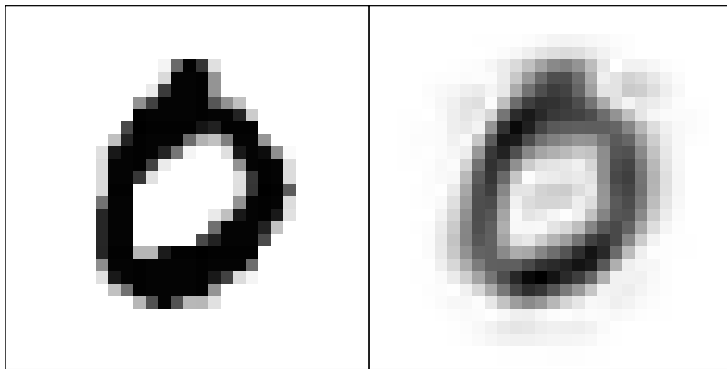- Consumption reflects some social structures

# Silver lining

In many settings data is often much more low-dimensional and not uniformly spread!

- Images have structures

- Biological systems are strongly regulated

- Consumption reflects some social structures

- etc.

# PCA on MNIST

# PCA on MNIST

Reconstructing the MNIST digits from the first 25 PCs

Overall, we need to be careful when analyzing high-dimensional data

Overall, we need to be careful when analyzing high-dimensional data

Methods that can identify relevant structures or incorporate the knowledge that these spaces are often sparse generally outperform traditional methods.

# Reference I

[Gir21]   Christophe Giraud. *Introduction to High-Dimensional Statistics*.
          Chapman and Hall/CRC, Aug. 2021.